

Les 3. Steekproefonderzoek

Onderzoekstechnieken

Jens Buysse Wim De Bruyn Wim Goedertier Bert Van Vreckem
AJ 2018-2019

**HO
GENT**

What's on the menu today?

Steekproefonderzoek

Kansverdeling van een steekproef

De Centrale Limietstelling

Van steekproef naar populatie

- Betrouwbaarheidsintervallen

- Betrouwbaarheidsinterval grote steekproef

- Kleine steekproef

- Betrouwbaarheidsinterval voor fractie

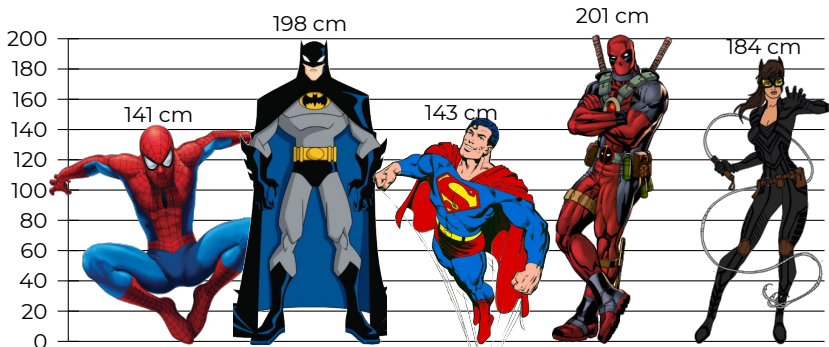
Steekproefonderzoek

USA Today has come out with a new survey. Apparently, three out of every four people make up 75% of the population

—David Letterman

Steekproef en Populatie

Herinner je onze superhelden



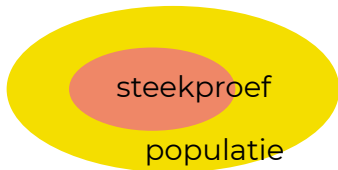
Steekproef en Populatie



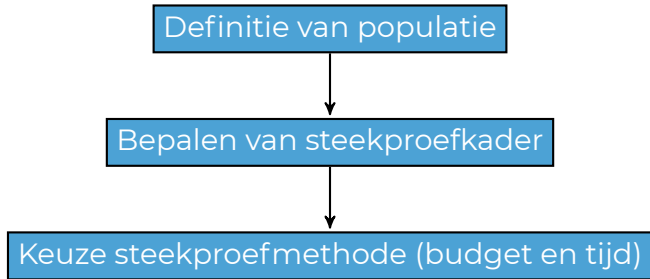
Steekproef en Populatie

De verzameling van alle objecten of personen waar men in geïnteresseerd is voor het onderzoek noemt men de **populatie**.

Wanneer men een subgroep uit een populatie gaat onderzoeken, dan noemen we die groep een **steekproef**.



Methode om tot een steekproef te komen



Gestratificeerd naar variabelen

Geslacht	Leeftijd				Totaal
	≤ 18	$]18, 25]$	$]25, 40]$	> 40	
Vrouw	500	1500	1000	250	3250
Man	400	1200	800	160	2560
Totaal	900	2700	1800	410	5810

Gestratificeerd naar variabelen

Geslacht	Leeftijd				Totaal
	≤ 18	$]18, 25]$	$]25, 40]$	> 40	
Vrouw	500	1500	1000	250	3250
Man	400	1200	800	160	2560
Totaal	900	2700	1800	410	5810

Geslacht	Leeftijd				Totaal
	≤ 18	$]18, 25]$	$]25, 40]$	> 40	
Vrouw	50	150	100	25	325
Man	40	120	80	16	256
Totaal	90	270	180	41	581

Hoe elementen voor een steekproef kiezen?

Aselecte (of random) steekproef : elk element uit de onderzoekspopulatie heeft een even grote kans om in de steekproef terecht te komen

Selecte steekproef : de elementen voor de steekproef worden *niet* random geselecteerd. Objecten die *gemakkelijk* kunnen verzameld worden, hebben meer kans om in de steekproef terecht te komen. (convenience sampling)



Mogelijke fouten

Metingen in een steekproef zullen typisch afwijken van de waarde in de hele populatie \Rightarrow Fouten!

- Toevallig \leftrightarrow Systematisch
- Steekproeffout \leftrightarrow Niet-steekproeffout

Steekproeffouten

- Toevallige steekproeffouten
 - Puur toeval

Steekproeffouten

- Toevallige steekproeffouten
 - Puur toeval
- Systematische steekproeffouten
 - Online enquête: mensen zonder internet worden uitgesloten
 - Straatenquête: enkel wie daar op dat moment loopt
 - Vrijwillige enquête: enkel geïnteresseerden doen mee

Niet-steekproeffouten

- Toevallige niet-steekproeffouten
 - Verkeerd aangekruiste antwoorden

Niet-steekproeffouten

- Toevallige niet-steekproeffouten
 - Verkeerd aangekruiste antwoorden
- Systematische niet-steekproeffouten
 - Slechte of niet geijkte meetapparatuur (slechte weegschaal)
 - Waarde kan beïnvloed worden door het feit dat je meet
 - Respondenten liegen (aantal sigaretten per dag)

Variantie/standaarddeviatie van een steekproef

Aangepaste formule:

$$s^2 = \frac{1}{n-1} \sum_{i=1}^n (\bar{x} - x_i)^2$$

De reden voor de wijziging kan je wiskundig bewijzen, maar we gaan het empirisch onderzoeken

R-script: `cursus/data/sample-variance.R`

Overzicht symbolen

	populatie	steekproef
aantal objecten	N	n
gemiddelde	μ	\bar{x}
variantie	$\sigma^2 = \frac{\sum(\mu - x_i)^2}{N}$	$s^2 = \frac{\sum(\bar{x} - x_i)^2}{n-1}$
standaarddeviatie	σ	s

Kansverdeling van een steekproef

Wat weten we nog van de kansrekening?

- Uitkomstenruimte
- Uitkomst
- Gebeurtenis
- Kansruimte



Kansverdeling voor één dobbelsteen

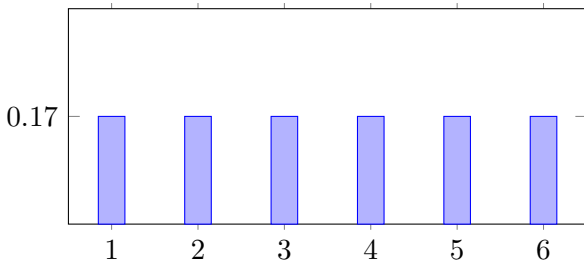
Wat is de kans om een aantal ogen te werpen met een dobbelsteen?

1	2	3	4	5	6

Kansverdeling voor één dobbelsteen

Wat is de kans om een aantal ogen te werpen met een dobbelsteen?

1	2	3	4	5	6
$\frac{1}{6}$	$\frac{1}{6}$	$\frac{1}{6}$	$\frac{1}{6}$	$\frac{1}{6}$	$\frac{1}{6}$

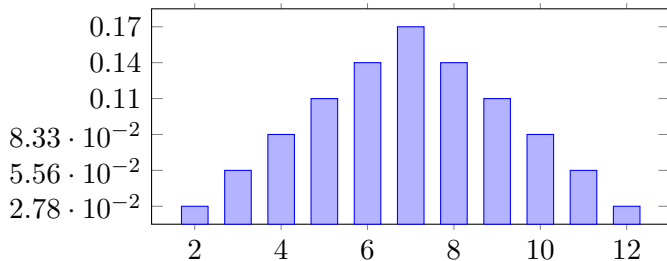


Kansverdeling voor twee dobbelstenen

2	3	4	5	6	7	8	9	10	11	12

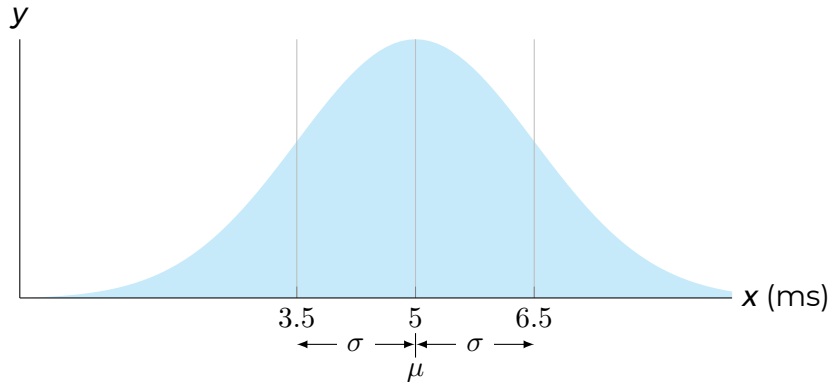
Kansverdeling voor twee dobbelstenen

2	3	4	5	6	7	8	9	10	11	12
$\frac{1}{36}$	$\frac{2}{36}$	$\frac{3}{36}$	$\frac{4}{36}$	$\frac{5}{36}$	$\frac{6}{36}$	$\frac{5}{36}$	$\frac{4}{36}$	$\frac{3}{36}$	$\frac{2}{36}$	$\frac{1}{36}$



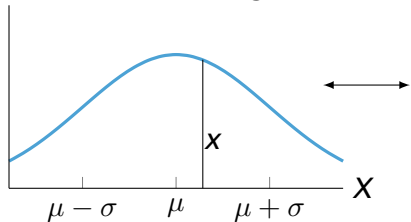
Continue kansverdeling

De reactiesnelheid x van Superman (in milliseconden) kun je als volgt weergegeven:



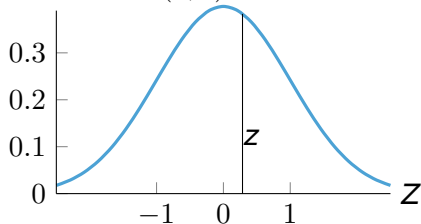
Standaardnormale verdeling

normaalverdeling $x \in X \sim \text{Nor}(\mu, \sigma)$



standaard normaalverdeling

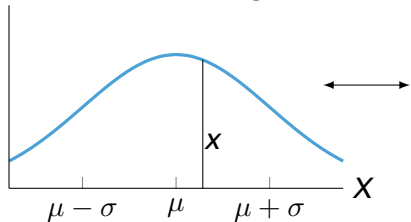
$z \in Z \sim \text{Nor}(0, 1)$



x en z hebben een vergelijkbare positie op de Gauss-curve.
Wat is het wiskundig verband tussen x en z ?

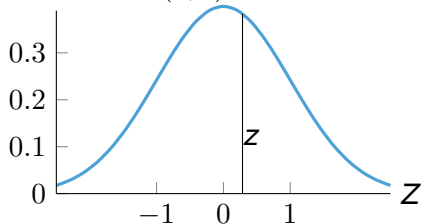
Standaardnormale verdeling

normaalverdeling $x \in X \sim \text{Nor}(\mu, \sigma)$



standaard normaalverdeling

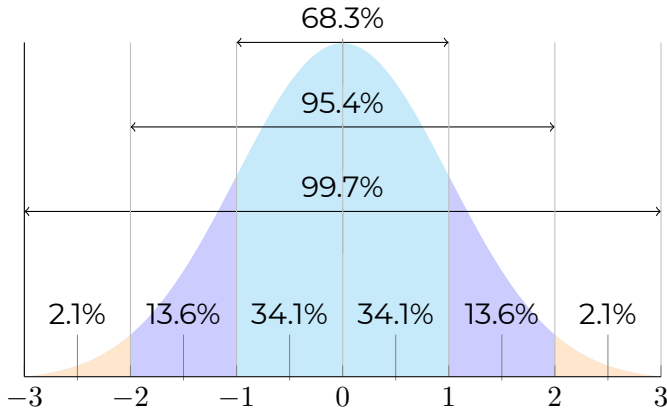
$z \in Z \sim \text{Nor}(0, 1)$



x en z hebben een vergelijkbare positie op de Gauss-curve.
Wat is het wiskundig verband tussen x en z ?

$$X = \mu + Z \cdot \sigma \quad \text{and} \quad Z = \frac{X - \mu}{\sigma}$$

Standaardnormale verdeling



Belangrijkste functies in R

Voor een normale verdeling met gemiddelde m en standaardafwijking s :

Functie	Betekenis
<code>pnorm(x, m, s)</code>	Linkerstaartkans, $P(X < x)$
<code>dnorm(x, m, s)</code>	Hoogte van de Gausscurve op punt x
<code>qnorm(p, m, s)</code>	Onder welke grens zal $p\%$ van de waarnemingen liggen?
<code>rnorm(n, m, s)</code>	Genereer n normaal verdeelde random getallen

Argumenten m en s weglaten geeft waarden voor de **standaard** normaalverdeling: `pnorm(x)=pnorm(x,0,1)`

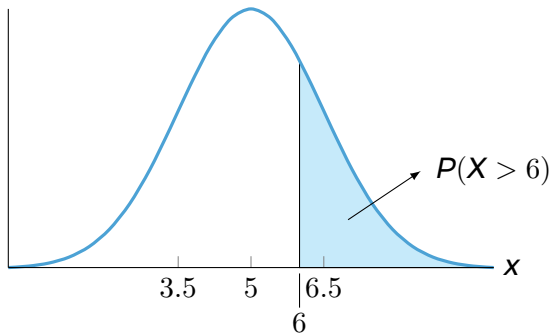
Kansen berekenen

Hoe groot is de kans dat...

...de reactiesnelheid van Superman meer dan 6 milliseconden is?

Wiskundige notatie:

$$P(X > 6) = ? \quad \text{met } X \sim \text{Nor}(\mu = 5, \sigma = 1,5)$$



Kansen berekenen: z-tabel

$P(X > 6) = ?$ met $X \sim \text{Nor}(\mu = 5, \sigma = 1,5)$

(Oude) berekeningsmethode met z-tabel, vb.

<http://sixsigmastudyguide.com/wp-content/uploads/2014/04/z-table.jpg>

Kansen berekenen: z-tabel

$$P(X > 6) = ? \quad \text{met } X \sim \text{Nor}(\mu = 5, \sigma = 1,5)$$

(Oude) berekeningsmethode met z-tabel, vb.

<http://sixsigmastudyguide.com/wp-content/uploads/2014/04/z-table.jpg>

1. zet om naar z-score

$$z = \frac{6-5}{1,5} = 0,667 \text{ dus } P(X > 6) = P(Z > 0,667)$$

2. zet eerst om naar een **linker**staartkans

- regel van 100% kans: $P(Z > 0,667) = 1 - P(Z < 0,667)$
- of symmetrie-regel: $P(Z > 0,667) = P(Z < -0,667)$

3. zoek op in z-tabel

Kansen berekenen: met R

$P(X > 6) = ?$ met $X \sim \text{Nor}(\mu = 5, \sigma = 1, 5)$

Rechtstreekse berekening met R (of rekenmachine)

Kansen berekenen: met R

$P(X > 6) = ?$ met $X \sim \text{Nor}(\mu = 5, \sigma = 1, 5)$

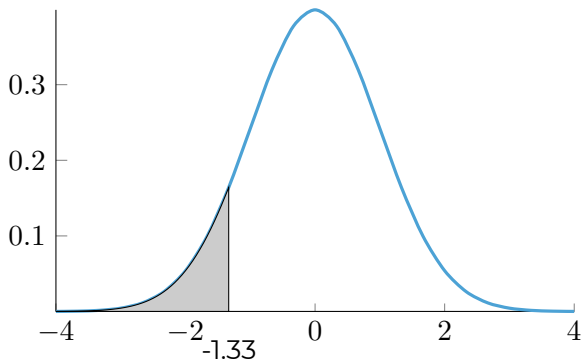
Rechtstreekse berekening met R (of rekenmachine)

- zet eerst om naar een **linker**staartkans:
met regel van 100% kans: $P(X > 6) = 1 - P(X < 6)$
- bereken met R: $1 - P(X < 6) = 1 - \text{pnorm}(6, 5, 1.5)$

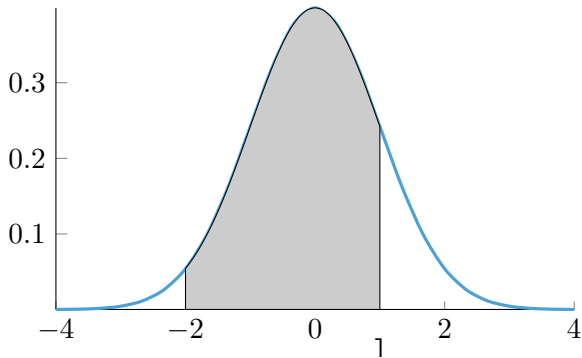
Voorbeelden

1. Hoe groot is de kans dat de reactiesnelheid van Superman minder dan 4 ms is?
2. Hoe groot is de kans dat hij in minder dan 7 ms reageert?
3. Hoe groot is de kans dat Superman in minder dan 3 ms reageert?
4. Hoe groot is de kans dat hij reageert tussen de 2 en de 6,5 ms?
5. Binnen welke tijd ligt 80% van zijn reactiesnelheid?

Vraag 3: $P(X < 3)$

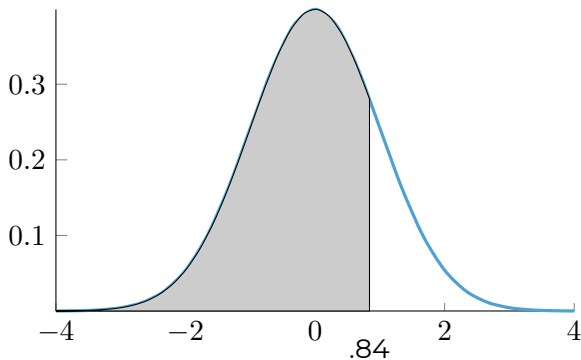


Vraag 4: $P(2 < X < 6,5)$



Vraag 5

Voor welke x is $P(X < x) = 80\%$?



De Centrale Limietstelling

De centrale limietstelling

Als de steekproefomvang voldoende groot is, dan kan de kansverdeling van het steekproefgemiddelde benaderd worden met een normale verdeling. Dit geldt ongeacht de vorm van de kansverdeling van de onderliggende populatie



- 1 test
- 25 tests
- 100 tests



Demo: <https://students.brown.edu/seeing-theory/probability-distributions/index.html>

De centrale limietstelling

Beschouw een aselechte steekproef van n waarnemingen die uit een populatie met een willekeurige distributie met verwachtingswaarde μ en standaardafwijking σ . Als n groot genoeg is zal de kansverdeling van \bar{x} een normale verdeling benaderen met gemiddelde $\mu_{\bar{x}} = \mu$ en standaardafwijking $\sigma_{\bar{x}} = \frac{\sigma}{\sqrt{n}}$.

Hoe groter de steekproef is, des te beter zal de kansverdeling van \bar{x} overeenkomen met een normaalverdeling.

Van steekproef naar populatie

Puntschatter

Een **puntschatter** voor een populatieparameter is een regel of een formule die ons zegt hoe we uit de steekproef een getal moeten berekenen om de populatieparameter te schatten.

Betrouwbaarheidsinterval

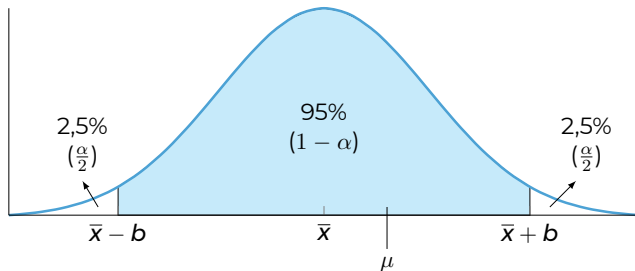
Een **betrouwbaarheidsinterval** is een regel of een formule die ons zegt hoe we uit de steekproef een interval kunnen berekenen dat de waarde van de parameter met een zeker **betrouwbaarheidsniveau** bevat.

Grote steekproef

Gegeven een steekproef met gemiddelde \bar{x} .

We zoeken nu een interval $[\bar{x} - b, \bar{x} + b]$ waarvan we met een betrouwbaarheidsniveau $(1 - \alpha)$ van bvb. 95% kunnen zeggen dat μ erbinnen ligt.

$$P(\bar{x} - b < \mu < \bar{x} + b) = 1 - \alpha = 0,95$$

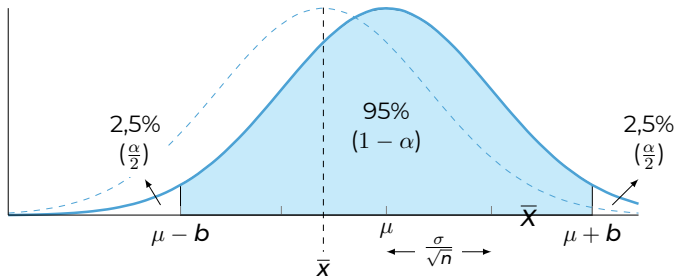


Grote steekproef

Dankzij de centrale limietstelling weten we dat $\bar{x} \in \bar{X} \sim \text{Nor}(\mu, \frac{\sigma}{\sqrt{n}})$

En omwille van de symmetrie weten we dat

$$P(\bar{x} - b < \mu < \bar{x} + b) = P(\mu - b < \bar{x} < \mu + b)$$



Grote steekproef

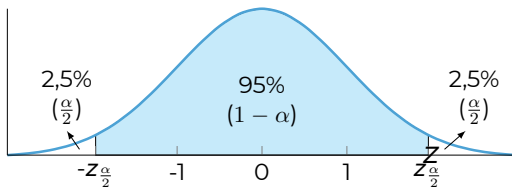
We bepalen nu de z-score voor \bar{x} : $z = \frac{\bar{x} - \mu}{\frac{\sigma}{\sqrt{n}}}$

We zoeken (in een tabel of rekentoestel) de waarde $z_{\frac{\alpha}{2}}$ waarbij:

$$P\left(-z_{\frac{\alpha}{2}} < z < z_{\frac{\alpha}{2}}\right) = 1 - \alpha = 0.95$$

$$P\left(z < z_{\frac{\alpha}{2}}\right) = 1 - \frac{\alpha}{2} = 0.975$$

$$z_{\frac{\alpha}{2}} = \text{qnorm}(0.975) \approx 1.96$$



Grote steekproef

We krijgen nu:

$$P\left(-1,96 < \frac{\bar{x} - \mu}{\frac{\sigma}{\sqrt{n}}} < 1,96\right) = 0,95$$

$$P\left(\mu - 1,96 \frac{\sigma}{\sqrt{n}} < \bar{x} < \mu + 1,96 \frac{\sigma}{\sqrt{n}}\right) = 0,95$$

Wegens symmetrie kunnen we μ en \bar{x} verwisselen:

$$P\left(\bar{x} - 1,96 \frac{\sigma}{\sqrt{n}} < \mu < \bar{x} + 1,96 \frac{\sigma}{\sqrt{n}}\right) = 0,95$$

We kunnen nu met 95% zekerheid zeggen dat:

$$\mu \in \left[\bar{x} - 1,96 \cdot \frac{\sigma}{\sqrt{n}}, \bar{x} + 1,96 \cdot \frac{\sigma}{\sqrt{n}} \right]$$

(in de praktijk gebruiken we $s_{steekpr}$ als schatting voor $\sigma_{populatie}$)

Kleine steekproef

Bij een ***kleine*** steekproef ($n < 30$) is de centrale limietstelling **niet** meer geldig.

Er geldt wél:

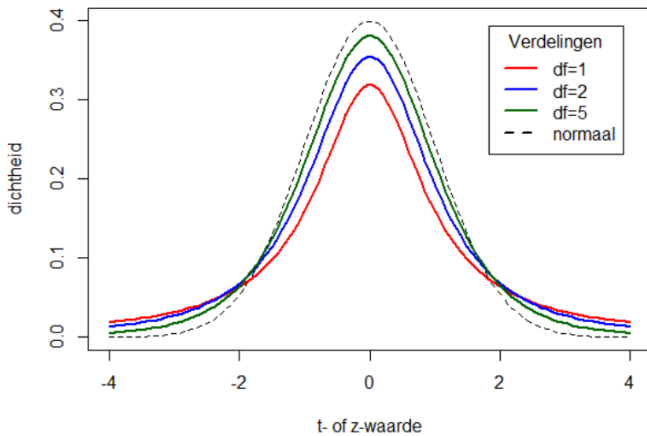
Als een populatie X normaal verdeeld is ($X \sim \text{Nor}(\mu, \sigma)$) en je neemt een ***kleine*** steekproef met gemiddelde \bar{x} en standaardafwijking s , dan is

$$t = \frac{\bar{x} - \mu}{\frac{s}{\sqrt{n}}}$$

verdeeld volgens een t-verdeling met $n - 1$ vrijheidsgraden (degrees of freedom)

De Student t -verdeling

Vergelijking van Student-t verdelingen



Student t -verdeling in R

Voor een t -verdeling met df vrijheidsgraden:
(df = degrees of freedom)

Functie	Betekenis
<code>pt(x, df)</code>	Linkerstaartkans, $P(X < x)$
<code>dt(x, df)</code>	Hoogte van de curve op punt x
<code>qt(p, df)</code>	Onder welke grens zal $p\%$ van de waarnemingen liggen?
<code>rt(n, df)</code>	Genereer n random getallen volgens deze verdeling

Betrouwbaarheidsinterval voor *kleine* steekproef

Om een betrouwbaarheidsinterval voor het gemiddelde μ van een populatie te bepalen op basis van een *kleine* steekproef, zoeken we eerst $t_{\frac{\alpha}{2}, df}$.

Bij een betrouwbaarheidsniveau van 95% is $\frac{\alpha}{2} = 0,025$

Veronderstel $n = 5$ (dus $df=4$), dan is

$$t_{\frac{\alpha}{2}, df} = qt(0.975, 4) = 2.776$$

We kunnen dan met 95% zekerheid zeggen dat:

$$\mu \in \left[\bar{X} - t_{\frac{\alpha}{2}, df} \cdot \frac{s}{\sqrt{n}}, \bar{X} + t_{\frac{\alpha}{2}, df} \cdot \frac{s}{\sqrt{n}} \right]$$

Betrouwbaarheidsinterval voor fractie

$$\bar{p} = \frac{\text{aantal successen}}{n}$$

- Verwachting van kansverdeling van \bar{p} is p .
- De standaardafwijking van kansverdeling \bar{p} is $\sqrt{\frac{pq}{n}}$
- Voor grote steekproeven is \bar{p} bij benadering normaal verdeeld.

Aangezien \bar{p} een steekproefgemiddelde is van het aantal successen, stelt dit ons in staat een betrouwbaarheidsinterval te berekenen analoog als die voor de intervalschatting van μ voor grote steekproeven.

$$\bar{p} \pm z_{\frac{\alpha}{2}} \cdot \sqrt{\frac{\bar{p}\bar{q}}{n}}$$

45/45 met $\bar{p} = \frac{x}{n}$ en $\bar{q} = 1 - \bar{p}$